## CONFIDENCE LIMITS FOR BIOLOGICAL ASSAYS
### C. I. Bliss

Connecticut Agricultural Experiment Station and Yale University

Most biological measurements of potency are based directly or indirectly upon a dosage-response curve (4). The experimenter selects and controls each dose so that its value is known within relatively narrow limits. The dose in turn determines the response, not precisely, but subject to the inevitable biological variation or sampling error. Usually the response is expressed in units $(y)$ which can be plotted as a straight line against the logarithm of the dose $(x)$. The line which best describes the relation of $y$ to $x$ is then computed by the method of least squares. As a straight line, it is defined by two independent statistics which determine its position $(a)$ and its slope $(b)$. The response corresponding to any log-dose can be calculated from the equation of the curve. While this computed response is not the true or population value, we can obtain upper and lower limits which will bracket the true value in a selected proportion $(P)$ of similar experiments. These are known either as fiducial or as confidence limits; in the present case both terms are used. When plotted over a suitable range they form two hyperbolae which at any given $x$ are equidistant above and below the fitted line in terms of $y$ (24). They are illustrated in Figure 1 by an experiment with tincture of digitalis on the relation between place of injection and the production of systolic standstill in the frog heart (20).

When measuring potency, a given response $(Y)$ or a difference in response is converted to units of log-dose $(X)$ by means of the dosage-response curve. Thus in Figure 1 the value of $X$ at the intersection of the dosage-effect curve with 5 probits or 50 percent is the log-LD50. If the susceptibility of frogs to digitalis were stable in different laboratories, the intersection of the limiting hyperbolae with the ordinate for 5 probits would define confidence limits which would bracket the true LD50 in 19 out of 20 similar experiments. In the present case these limits are not spaced equally to the right and left of the most probable value. Similar limits of the log-LD5 (Figure 1) are spaced even more unequally, but with the large and small intervals reversed. This sort of inequality is neglected in computing limits from the standard error of the log-LD50 or of the log-ratio of potencies $(s_M)$. In the large sample approximation the standard error is multiplied by the value of $t$ at the required level of $P$ and the product is added to or subtracted from the most probable value. Not only are two intervals averaged which may be quite unequal but the average may be underestimated. In critical cases, this may lead to the wrong conclusion.

The true confidence limits could be obtained graphically by interpolation from hyperbolae such as those in Figure 1 but it is easier to compute them directly. They have been

used frequently in English papers on biological assay but seldom in this country. Apparently the true limits were published first in 1935 (1) as applied to the log-dose for an assigned response, such as for the log-LD50. Later, the interpretation of confidence limits and their derivation were discussed by Eisenhart (7), who gave equations suitable for biological assays. The derivation of the fiducial or confidence limits for biological assays were reported also by Fieller (8) in 1940 and by Irwin (16) in 1943. More recently these equations have appeared in a number of papers, in varying symbolism. They will be converted here to a uniform symbolism for application to the various types of biological assays.

*The terms requiring confidence limits.* We will assume that the response has been converted to values $(y)$ which can be plotted linearly against the log-dose $(x)$ over an adequate dosage range, so that they agree with the equation for a straight line

$$Y = a + b(X - \bar{x}), \qquad \ldots (1)$$

where $a$ is equal numerically to the observed mean response $\bar{y}$, $\bar{x}$ is the observed mean log-dose and the slope is given by the equation

$$b = \frac{[xy]}{[x^2]} = \frac{S(xy) - \bar{y}S(x)}{S(x^2) - \bar{x}S(x)} \qquad \ldots (2)$$

In some cases, such as with the all-or-none response, weights $(w)$ are used in computing both the means and the slope and in this case $b = [wxy] / [wx^2]$. By rearrangement of equation (1), the most probable log-dose at any given response $Y$ is

$$X = \bar{x} + (Y - \bar{y}) / b \qquad \ldots (3)$$

If a single value is required, such as for the log-LD50, it is determined by equation (3) or its equivalent and *not* as the mid-range of the exact limits.

The experimenter is usually more interested in relative than in absolute values, typically in the potency of an unknown preparation relative to a standard. Two dosage-response curves may be determined such as are illustrated by the two curves in Figure 1, which compare two routes of administration for tincture of digitalis rather than two drugs. In the absence of a qualitative difference in their activity, the dosage-response curves for "Standard" (lymph sac) and "Unknown" (muscle) should be parallel within the sampling error. This assumption is confirmed by statistical test (1, 5, 6, 8, 10, 20). Then the slope of the best-fitting pair of parallel lines is computed from the sums of the numerators and of the denominators of the component curves as

$$b_c = S [xy] / S [x^2]. \qquad \ldots (2a)$$

When the data of many similar experiments are available, all agreeing in slope within the sampling error, the combined slope may be used in analyzing the individual tests (3, 8, 15). This reduces the curvature in the limiting hyperbolae. Those plotted in Figure 1, for example, are based upon the standard error of the slope for the combined series; they would have shown a greater curvature if each pair of limits had been computed independently.

The logarithm of the potency of the unknown relative to that of the standard is given by the equation

$$M = \bar{x}_s - \bar{x}_u - \frac{\bar{y}_s - \bar{y}_u}{b_c} = \bar{x}_s - \bar{x}_u + \frac{\bar{y}_u - \bar{y}_s}{b_c} \qquad \ldots (4)$$

where $\bar{x}_s$ and $\bar{y}_s$, $\bar{x}_u$ and $\bar{y}_u$ are the mean log-doses and responses for standard and unknown respectively. This is the form of $M$ in most common use, although Finney (10) has interchanged the subscripts in equation (4) so that his $M$ is equivalent to $-M$ in the above notation. Apparently the change in sign has led to the discrepancy between the equations for the mid-

interval of the exact confidence limits given by Gridgeman (13, 14) and those in a later section of the present paper. The most probable single estimate of potency for a given experiment is that determined from equation (4). The unknown is assigned an assumed potency and equivalent doses of Standard and Unknown are administered. When the doses of Unknown are expressed in assumed units, the mean log-dose of Standard and Unknown are often equal, so that $\overline{x}_s - \overline{x}_u = 0$. The log-ratio of potencies may then be written as

$$M' = (\overline{y}_u - \overline{y}_s) \ / \ b_c. \qquad \qquad \ldots \ (5)$$

These are the values for which confidence limits are needed.

*Approximate vs. exact confidence limits.* The approximate method for establishing confidence limits is to determine the standard error for $X$ or $M$, multiply the standard error by Student's $t$ as read from a suitable table (11) at the required level of probability $(P)$, and add the product to and subtract it from the most probable value, giving $X \pm t s_x$ and $M \pm t s_M$ (2, 3, 5, 20, 22). These limits have the advantage of being somewhat easier to compute, more familiar, and in many cases a very close approximation to the exact limits. In determining the ED50, for example, the observations may fall more or less equally above and below the required value, so that $\overline{y}$ approaches 5 probits where the approximate and exact limits nearly coincide. In assays of relative potency, the difference $(\overline{y}_s - \overline{y}_u)$ may be so small that there is little gain from using the exact limits. This is shown in Figure 1 by drawing the line marked "M" midway between $\overline{y}_s$ and $\overline{y}_u$. At this level of response the confidence intervals to the right and left of the respective curves are nearly equal. When there is information from many separate, consistent estimates of the slope, the variance of the pooled value is small and the approximate limits agree well with the exact ones. In many cases, therefore, little is gained by substituting exact confidence limits for their approximate values.

The arguments in favor of the exact limits, however, are still cogent. With a single dosage-response curve or with a single assay based upon relatively few animals, the slope of the dosage-response curve may have a large sampling error. In some cases $b$ may not differ significantly from 0 at the required level of significance especially where $P \leqslant 0.01$ (8). As Irwin has pointed out (16, 17), the confidence limits are then indeterminate. This accords with common sense; an assay in which the known differences in potency cannot be detected is hardly adequate for measuring the unknown differences between materials. When the dose must be predicted for a response $(Y)$ which differs considerably from $\overline{y}$, the exact limits are preferred. Thus in Figure 1 the confidence interval at LD5 for injection in the lymph sac differs more from the approximate value than that for LD50 and the same would be true in estimating 95 or 99 percent kill from toxicity tests of insecticides. In determining the safety of a drug the chemotherapeutic index may take the slope into account by comparing the ED99 with the LD1 (12) or the ED95 with the LD5 (23), where ED indicates the effective dose at a given percent response just as LD stands for the lethal dose. Here again the error of the slope is important. The use of the approximate limit would both underestimate the width of the confidence interval and bias its position (17). It must be emphasized again, however, that the midpoints of the exact confidence intervals are not the most probable values. These are given by equations (3), (4) and (5) and are the midpoints of the approximate intervals.

*Limits of the log-dose for a given graded response.* The first requirement for a valid experiment is that the estimated slope differ significantly from 0 at the confidence level which is specified. This can be tested most conveniently (8) by computing the difference

$$B^2 \ - \ s^2 t^2. \qquad \qquad \ldots \ (6)$$

$B^2$ measures the variation in $y$ accounted for by the slope of the dosage-response curve where

$$B^2 = \frac{[xy]^2}{[x^2]} \text{, or } B^2 = \frac{S^2[xy]}{S[x^2]} . \qquad \ldots (7)$$

The second form is used when the data from several dosage-response curves are pooled in computing $b_c$. The value of $t$ is read from a table of Student's distribution (11) at the appropriate level of $P$ with $n$ equal to the degrees of freedom of the variance for error, $s^2$. If the difference $(B^2 - s^2t^2)$ is positive, the slope differs significantly from 0 and confidence limits can be determined.

The next step is to compute a correction term $C^2$ ($C_p$ in Fieller's notation (8) ) which indicates the extent of the discrepancy between approximate and exact limits. Whenever approximate limits are used $C^2$ is assumed to equal 1; exact limits are calculated with the observed value of $C^2$. It is determined as

$$C^2 = B^2 / (B^2 - s^2t^2). \qquad \ldots (8)$$

We may now compute the exact confidence limits which will enclose the true or population value of $X$ at a given response $Y$ in the proportion $P$ of similar experiments. Usually the level of $P = 0.95$ or $0.99$ is used, which means that in $0.05$ or $0.01$ of similar experiments the true or population value will fall outside the computed limits. The value of $t$ is read from the table at $P = 0.05$ or $P = 0.01$, not at $0.95$ or $0.99$. The exact limits are then given by the equation

$$X_L = \overline{x} + C^2(Y - \overline{y}) / b \pm t\lambda C \sqrt{1/N + (Y - \overline{y})^2 / (B^2 - s^2t^2)}, \qquad \ldots (9)$$

where $N$ is the number of observations in computing the regression line and $\lambda = s/b$, the other terms having their former meaning. This is equivalent to Equation 30 in Fieller's paper (8). If weighting coefficients have been used in computing the curve, $N$ is replaced by $S(w)$. The derivation of the equation is given by Irwin (16).

The response indicated by $Y$ in equation (9) is assumed to be known without error. It has been applied to the assay of insulin by Fieller (8). Six different dosage-levels of the Unknown were compared in separate cross-over tests with a single level of Reference Standard, fitting a straight line to the differences in response $(y)$. The log-potency of the Unknown was then equal to $X$ at $Y = 0$. However, these confidence limits do not represent the predictive value of the equation in meeting a tolerance requirement. The problem has been discussed at length by Eisenhart (7). A third term, the digit 1, must be added to the expression under the radical in equation (9) and $t$ must be a "one sided $t$", read at $P = 0.10$, if one wishes odds of $0.95$ that the observed result will be larger than some specified value.

The exact differ from the approximate limits only in the introduction $C$ (or $C^2$) and of the product $s^2t^2$ under the radical. When the slope accounts for a large fraction of the total variation, so that it is highly significant, $C^2$ is but little larger than 1 and the approximate limits are usually adequate.

*Limits of the log-dose for a given all-or-none effect.* With an all-or-none response the dosage-effect curve is based upon the percentage of positive reactions in independent groups tested at different doses. Originally this curve has an asymmetrical sigmoid shape but in most cases it can be changed to a straight line by changing percentages to probits and plotting probits against the log-dose. For the most accurate results, the line is computed by the method of maximum likelihood with corrected probits and weights (2). The calculation enables one to compare by the $\chi^2$ test the variance based upon the discrepancy between the observed values and the regression line with that expected from the binomial distribution. If $\chi^2$ indicates satisfactory homogeneity, the value of $t$ is that for $n = \infty$, that is, the normal deviate. Since the response is in units of standard deviations, $s^2 = 1$, and a positive $B^2 - t^2$

indicates that the slope is significant. The value of $C^2$ becomes $C^2 = B^2 / (B^2 - t^2)$.

The exact confidence limits for a given response $(Y)$ in probits are

$$X_L = \bar{x} + \lambda C^2(Y - \bar{y}) \pm t\lambda C\sqrt{1/S(w) + (Y - \bar{y})^2/(B^2 - t^2)}, \qquad \ldots \text{(10)}$$

where $S(w)$ is the sum of the weights used in computing the curve. Since $s^2 = s = 1$, $\lambda$ equals $1/b$. For the log-LD50, the equation is solved with $Y = 5$. Occasionally $\chi^2$ will indicate a greater scatter of the observations about the fittted line than would be expected by chance, although inspection of the plotted points shows that no simple curve fits better than a straight line. Under these circumstances the term following the $\pm$ is multiplied by $\sqrt{\chi^2/n}$ and $t$ is that for the degrees of freedom $(n)$ of $\chi^2$.

Irwin (16) has compared the approximate with the exact limits for eleven determinations of the median fertility dose of vitamin E. As expected, the limits agreed more closely at $P = 0.95$ that at $P = 0.99$ but even at $P = 0.95$ the two estimates differed markedly in four of the eleven tests. When the example in Figure 1 for injection into the lymph sac was computed alone (not with the combined slope) the LD50 = 8.56 cc/kg with approximate limits of 6.95 and 10.52 and exact limits of 7.31 and 27.57 by equation 10. When the combined slope was used instead, the LD50 = 8.37 cc/kg with approximate and exact limits of 7.19 and 9.74, and 7.31 and 10.59 respectively. The advantage of pooling the information on slope as in Figure 1 is obvious. In Irwin's paper (16) the limits are stated in percentages of the expected LD50.

*Limits for log-ratios of potency.* In measuring the potency of an Unknown relative to a Standard, the experimenter usually computes two dosage-response curves; one for each material. If these curves are parallel when dosage is expressed in logarithms, the potency measured at one level of response is the same as that at any other, so that the Unknown "stays standardized". In some cases, however, the log-dose at a single level, such as the log-LD50, is computed separately from each curve and the potency determined from the two LD50's. This does not bias the estimate of relative potency but it can increase the estimate of error. In Figure 1, for example, both means fall well under 5 probits, so that the confidence limits at $Y = 5$ are wider than at a point mid-way between $\bar{y}_s$ and $\bar{y}_u$. Hence the confidence limits for $M$ are calculated so as to minimize the effect of errors in slope. Computing $M$ directly has the further advantage of pooling the data both on the slope and on the variance, so that these are estimated with greater reliability.

For an assay based upon a graded response the confidence limits as given by the British Standards Method for Vitamin $D_2$ (6) and by Irwin (16) are

$$X_L = \bar{x}_s - \bar{x}_u - C^2(\bar{y}_s - \bar{y}_u)/b_c \pm t\lambda C\sqrt{1/N_s + 1/N_u + (\bar{y}_s - \bar{y}_u)^2/(B^2 - s^2t^2)}, \quad \text{. (11)}$$

where $C^2$ is computed by equation (8) with the combined values of $B^2$ and $s^2$, and $t$ is that for the degrees of freedom of $s^2$. This equation differs from the approximate form,

$$X_L = M \pm ts_M \qquad \dots (12)$$

only in the introduction of $C^2$, $C$ and under the radical $s^2t^2$. An example of the calculation is given in detail with a somewhat different nomenclature by the British Standards Method for the biological assay of vitamin $D_3$ with chicks (6). The exact and approximate limits have been compared in a long series of vitamin A assays by Irwin (17). In a considerable number of cases the exact method gave indeterminate confidence limits, the approximate method leading to finite but unreliable results. When many assays are available, the slope and variance may be sufficiently homogeneous to warrant combining the entire evidence. Frequently this avoids the instability due to the inadequacy of the slope and the variance in any component assay.

The confidence limits based upon an all-or-none or quantal reaction are very similar to those in Equation (11) for a graded response. They are equal to

$$X_L = \bar{x}_s - \bar{x}_u - \lambda C^2(\bar{y}_s - \bar{y}_u) \pm t\lambda C \sqrt{1/S \ (w_s) + 1/S(w_u) + (\bar{y}_s - \bar{y}_u)^2 / (B^2 - t^2)}, \qquad \dots (13)$$

where $C^2$ is computed from equation (8) with $B^2$ based on the pooled slope. In a well-designed assay $\bar{y}_s$ differs but little from $\bar{y}_u$, so that the approximate limits are often adequate. Irwin (18) has computed the exact and approximate limits for a long series of insulin assays by the mouse convulsion test. In many cases they were run with two doses of the unknown and one dose of the standard so that $\bar{y}_s$ was replaced by $y_s$ and $S(w_s)$ by $w_s$. Although the slope varied significantly between assays, he averaged the slopes and enlarged the error by a factor equivalent to $\sqrt{\chi^2/n}$.

*Assays designed factorially*. It is often convenient to use a balanced factorial design for biological assays (5). These have an equal number of observations on the Standard and on the Unknown and on each dose, and equal intervals between successive log-doses. In terms of the amount of the Unknown which is assumed to equal one unit of Standard, the log-ratio of potencies may be computed as

$$M' = ikT_1/T_2, \qquad \dots (14)$$

where $k = 1$, $4/3$ and $5$ for assays with 2, 3 and 4 doses respectively of both Standard and Unknown and $i$ is the interval between successive log-doses. $T_1$ and $T_2$ are the sums of the products $S(xY_p)$ obtained by multiplying the totals $(Y_p)$ for each dose of each preparation by the factorial coefficients for differences between drugs, giving $T_1$, and by those for the combined slope, giving $T_2$.

The confidence limits for $M'$ are

$$X_L = C^2M' \pm t\lambda C \sqrt{4[1 + D^2/(B^2 - s^2t^2)]/N'S(x_1^2)} \qquad \dots (15)$$

where $N'$ is the number of observations for each dose of Standard and Unknown or the number of groups or "blocks", $S(x_1^2)$ is the sum of the squares of the coefficients used in computing $T_1$, $D^2$ is the variance due to the difference in response between the Standard and the Unknown or $D^2 = T_1^2/N'S(x_1^2)$, and the other symbols have the same meaning as before. For approximate limits, $X_{L'}$, $C$ is assumed to equal 1 and $s^2t^2$ is omitted from equation (15).

*Two-dose factorial assays*. The most generally useful factorial assays require two doses of the Standard and two doses of the Unknown, which are administered at random to $N'$ homo-

geneous groups or "blocks", each of four reactions. Several papers have described methods for their analysis (3, 10, 14, 22). While differing in nomenclature, they agree in the calculation of the relative potency of the Unknown by equation (14) as $M' = iT_1/T_2$. In the shortened calculation (3), $T_1$ is computed as the total of the $y_1$'s and $T_2$ as the total of the $y_2$'s, the $y$'s being determined separately for each complete group of four reactions. For computing $s^2$, information is also available from the third independent difference in a group of four which is designated as $y_3$.

The confidence limits may be written in a similar form for both approximate and exact values. Both involve a different correction term from that used before which will be defined as

$$c^2 = 4s^2t^2N', \qquad \ldots (16)$$

where $s^2$ is the variance for all components in the error term; i.e., the variation in $y_1$, in $y_2$ and in $y_3$, and $t^2$ is based upon the degrees of freedom of $s^2$. The approximate limits are

$$X_L{}^1 = i(T_1T_2 \pm c\sqrt{T_1{}^2 + T_2{}^2})/T_2{}^2. \qquad \ldots (17)$$

The exact limits are

$$X_L = i(T_1T_2 \pm c\sqrt{T_1{}^2 + T_2{}^2 - c^2}) / (T_2{}^2 - c^2). \qquad \ldots (18)$$

This exact equation was given first by Schild (22) based upon a memorandum from Irwin, and later by Finney (10), whose fiducial limits are for $-M$ instead of $M$. Gridgeman (14) gives the same range on each side of the mid-interval but apparently does not have the correct mid-interval.

Information can often be pooled from a number of two-dose assays, so that more data are available on the slope than on the difference between Standard and Unknown (14). Occasionally the three components of $s^2$, representing the variation in $y_1$, in $y_2$ and in $y_3$, are not equal. Thus in the assay of penicillin, Knudsen (19) has reported for one laboratory a significantly greater variance for the slope (in $y_2$) than for the difference between drugs (in $y_1$). When the biological material occurs naturally in pairs, one component may be confounded by the experimental design with the difference between the pairs, so that it is measured with less precision than the others. This procedure was followed by Price and Spencer (21) in the assay of virus activity. Equivalent concentrations of Standard and Unknown were applied to opposite halves of the same leaf and different concentrations to different leaves, so that the variance in the slope was larger than that for the comparison of Standard and Unknown at a given dosage level. A preferable design is that reported by Fieller (9), in which the test for parallelism ($y_3$) is confounded with differences between rabbits in the cross-over assay for insulin. For these and other reasons it may be necessary to separate the components which are pooled in equation (16).

The problem may be solved in a form comparable to equation (18) by the use of two "correction" terms similar to that defined by equation (16). The derivation is based upon the papers by Finney (10) and by Gridgeman (13) but with corrections which have been noted. It is assumed that $T_1$ is based upon $N_1$ groups each of four reactions with a variance of $s_1{}^2$ and a value of $t_1$ for the required $P$ with the degrees of freedom of $s_1{}^2$. A correction term is computed equal to

$$c_1{}^2 = 4s_1{}^2t_1{}^2N_1. \qquad \ldots (19a)$$

The total for slope, $T_2$, is computed from $N_2$ similar groups with a variance of $s_2{}^2$ and a value of $t_2$ at the same $P$ as for $t_1$ but determined for the degrees of freedom of $s_2{}^2$. These give a second correction term of
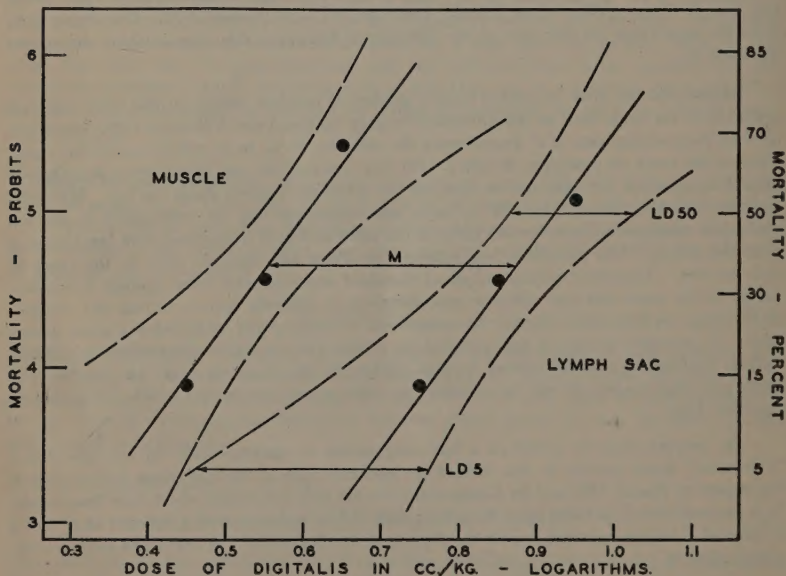
$$c_2{}^2 = 4s_2{}^2t_2{}^2N_2. \qquad \ldots (19b)$$

The confidence limits are then equal to
$$X_L = [T_1 T_2 \pm \sqrt{c_2^2 T_1^2 + c_1^2 (T_2^2 - c_2^2)}] \, iN_2 / N_1 (T_2^2 - c_2^2).$$
(20)

This reduces to equation (18) when $N_1 = N_2$ and $c_1 = c_2$.

In computing either exact or approximate confidence limits, it should be remembered that they concern only the sampling errors within individual assays. While the precision of an assay is no greater than that indicated by the limits, it may be considerably less as a result of factors which are constant within assays. An estimate of potency may be seriously impaired by factors such as failure to use an adequate technique in preparing the initial stock solutions or to guard against their deterioration. These depend upon the skill of the experimenter but a knowledge of the inherent precision of an estimate can be invaluable in finding and correcting these non-random sources of error.

*Summary.* A large-sample approximation may be involved when fiducial or confidence limits are computed from the standard errors of the terms used in biological assays. In many cases a correct answer requires the exact limits, which are not difficult to compute when their equations are reduced to a uniform symbolism. For experiments based upon a graded response the exact limits cover the log-dose for both a given and an observed response and for two-dose design and the confounded two-dose design. Equivalent limits based upon an all-or-none or quantal response include these for a required effect, such as the LD50 and for the log-ration of potencies.



Effect of route of administration upon the action of digitalis as measured by systolic standstill one hour after injection (20). Each point represents 15 frogs. The equations of the solid dosage-response lines are $Y = 4.609 + 6.840(X - 0.8655)$ for lymph sac ("Standard") and $Y = 4.714 + 6.840(X - 0.5613)$ for muscle ("Unknown"). The following additional terms are used in computing the confidence limits (broken lines) at $P = 0.05$: $1/S(w)_s = 0.04202$, $1/S(w)_u = 0.04049$, $1/S[wx^2] = 3.355$, $B^2 = 13.94$, $t = 1.960$, $C^2 = 1.380$.

## LITERATURE CITED

1. Bliss, C. I. The comparison of dosage-mortality data. Ann. Appl. Biol. **22**: 307-335 (1935).

2. Bliss, C. I. The determination of the dosage-mortality curve from small numbers. Quar. Jour. Pharm. and Pharmacol. **11**: 192-216 (1938).

3. Bliss, C. I. A simplified calculation of the potency of penicillin and other drugs assayed biologically with a graded response. Jour. Amer. Stat. Assn. **39**: 479-487 (1944).

4. Bliss, C. I. and McKeen Cattell. Biological Assay. Ann. Rev. Physiol. **5**: 479-539 (1943).

5. Bliss, C. I. and H. P. Marks. The biological assay of insulin. II. The estimation of drug potency from a graded response. Quar. Jour. Pharm. and Pharmacol. **12**: 182-205 (1939).

6. British Standard Method for the Biological Assay of Vitamin $D_3$ by the Chick Method British Standards Institution. Specification 911. (1940).

7. Eisenhart, C. The interpretation of certain regression methods and their use in biological and industrial research. Ann. Math. Stat. **10**: 162-186 (1939).

8. Fieller, E. C. The biological standardization of insulin. Suppl. Jour. Roy. Stat. Soc. **7**: 1-64 (1940).

9. Fieller, E. C. A fundamental formula in the statistics of biological assay, and some applications. Quar. Jour. Pharm. and Pharmacol. **17**: 108-123 (1944).

10. Finney, D. J. Mathematics of biological assay. Nature **153**: 284 (1944).

11. Fisher, R. A. and F. Yates. Statistical Tables for Biological, Agricultural and Medical Research. Second Edition. Oliver and Boyd, London (1943).

12 Foster, R. H. K. Standardization of the safety margin. Jour. Pharm. and Expt. Thera. **65**: 1-17 (1939).

13. Gridgeman, N. T. Mathematics of biological assay. Nature **153**: 461-462 (1944).

14. Gridgeman, N. T. The estimation of vitamin A. Lever Bros. and Unilever Ltd. London (1944).

15. Gridgeman, N. T. Special designs for vitamin D assays. Quar. Jour. Pharm. and Pharmacol. **18**: 15-23 (1945).

16. Irwin, J. O. On the calculation of the error of biological assay. Jour. Hygiene **43**: 121-128 (1943).

17. Irwin, J. O. A statistical examination of the accuracy of vitamin A assays. Jour. Hygiene **43**: 291-314 (1944).

18. Irwin, J. O. The error of the biological assay of insulin by the mouse convulsion test. Quar. Jour. Pharm. and Pharmacol. **16**: 352-362 (1943).

19. Knudsen, L. F. Control chart analysis of penicillin assays. Jour. Bact. August (1945).

20. Miller, L. C., C. I. Bliss and H. A. Braun. The assay of digitalis. Criteria for evaluating various methods using frogs. Jour. Amer. Pharm. Assn. **28**: 644-657 (1939).

21. Price, W. C. and E. L. Spencer. Accuracy of the local-lesion method for measuring virus activity. III. The standard deviation of the log-ratio of potencies as a measure of the accuracy of measurement. Amer. Jour. Bot. **30**: 720-735 (1943).

22. Schild, H. O. A method of conducting a biological assay on a preparation giving repeated graded responses illustrated by the estimation of histamine. Jour. Physiol. **101**: 115-130 (1942).

23. Whitlock, J. H. and C. I. Bliss. Bioassay technique for anthelmintics. Jour. Parasitol. **29**: 48-58 (1943).

24. Working, H. and H. Hotelling. The application of the theory of error to the interpretation of trends. Proc. Amer. Stat. Assn. p. 73-85 (1929).

# TRAINING IN STATISTICS AT THE UNIVERSITY OF MINNESOTA

Courses in statistics for undergraduates are offered in several of the colleges of the University. While basic statistical principles and techniques are taught in beginning courses in the different colleges, the illustrative problems are drawn from the field of interest of the students. Students who know, as undergraduates, that they expect to do graduate work in statis-tics are encouraged to secure a good founda-tion in mathematics during their undergrad-uate career.

Study at the graduate level in the University of Minnesota is organized through its Graduate School. Courses in statistics for graduate students are offered in seven different depart-ments. Such courses are available to students

who want to become professional statisticians and to the large group who want to learn some statistics to be used as a tool in their profession. Many advisors of graduate students of this latter class require that their advisees take one or more courses in statistics. There is a rapidly growing appreciation and use of statistical methods in the design of experiments and the interpretation of data. Whenever extensive data are collected statistics will be a useful tool in their interpretation. Scientific literature shows an increasing use of statistical methods for the interpretation of data. Unless the young research worker obtains training in statistics during his graduate student career, he will find that he is not only handicapped in designing proper experiments during his professional career but will find difficulty in understanding published literature in which modern statistical methods are used.

The Masters degree may be taken in the University of Minnesota with either a major or minor in statistics. The courses taken will usually comprise courses in mathematics and such courses in statistics as relate to the subject matter interests of the candidate. The courses from which the candidate may choose, with the approval of his advisor and a group committee of the Graduate School will be described under the plan for the Ph. D. degree.

A plan has been set up recently for granting the Doctor of Philosophy degree in statistics. For a major in statistics the candidate for the Ph.D. will be expected to present undergraduate preparation in mathematics up to and including differential and integral calculus. The course work should include a minimum of 50 quarter credits chosen from courses approved by a special faculty committee, but shall include two courses in Advanced Calculus and three courses in the Mathematical Theory of Statistics, or their equivalent. The minor program will include 21-24 quarter credits in a subject matter field in which the candidate expects to apply his statistical theory.

The candidate's program will be formulated with the aid of an advisor, selected from a special faculty committee, who is most closely associated with the subject matter field in which

the minor work is taken. This program will be evaluated by this special committee and then referred to the appropriate group committee of the Graduate School. A reading knowledge of two foreign languages is required.

Students whose major for the Ph.D. degree is one of the sciences and who look forward to research requiring statistical techniques may take a minor in statistics. The minimum number of credits will be 24, and the minor program will be worked out in consultation with an advisor selected from the special committee previously referred to.

Among the courses from which candidates for either a major or a minor in statistics may be selected are the following: differential equations, advanced calculus, mathematical theory of statistics, vector analysis, algebraic theory, theory of equations, topics in analysis and theory of functions in the Department of Mathematics; biometric principles, correlation analysis, statistical inferences, topics in biostatistics, life tables, research problems in biostatistics, seminar in biostatistics, and certain laboratory courses in Biostatistics; business statistics, correlation, index numbers, and senior topics in statistics in Business Administration; theory of statistics in Economics; advanced agricultural statistics in Agricultural Economics; methods in educational research, statistical methods in education, and problems in statistics in Educational Psychology; and a course in applied statistics in Plant Genetics. Some of the above courses are for a single quarter while others are for two or three quarters.

A special committee on the Ph.D. program in statistics has been set up by the graduate school. This committee consists of Dunham Jackson, Mathematics; P. O. Johnson, Education; B. D. Mudget, Business Administration; W. C. Waite, Agricultural Economics; A. E. Treloar, Biostatistics; and F. R. Immer, Plant Genetics. The candidate's program will be formulated with the aid of that advisor selected from this special committee who is most closely associated with the subject matter field in which the candidate expects to apply statistics.

# THE ECONOMICS OF SAMPLE SIZE APPLIED TO THE SCALING OF SAWLOGS*

## R. H. Blythe, Jr.

Division of Forest Economics, Forest Service, Department of Agriculture

On the national forests, payment for the sale of standing timber is based on measurement of the logs actually cut from the timber. The measurement is made by scaling the logs, which means measuring the length and the diameter at the small end. Standard tables showing the volume of lumber in board feet that may be cut from logs of different diameters and lengths are then used to find the total volume. Customary practice on national forests and on private sales has usually been to measure all of the logs in a sale, which sometimes is a very large number. Obviously, this is a situation where sampling has considerable possibilities.

Studies were conducted in most of the major timber producing regions of the country to establish the feasibility of sample scaling. These studies indicated that an efficient sampling system can be applied in the field.

*Variation in Log Sizes.* The variation in log sizes is, of course, the determining element in the possibility of sample scaling in general. Taking the country as a whole, variation in size is rather extreme, being from 0 to 6,500 board feet per log. However, the studies referred to above brought out that while the standard deviation of the volume of individual logs varied greatly from one part of the country to another and from one type of timber to another, the co-efficient of variation is remarkably constant over the country. It was found to range from .50 to .90. Thus, for purposes of practical application by persons untrained in statistical methods, rather simple tables can be made showing the number of logs that must be scaled in order to give a standard error value which administrative policies determine to be satisfactory both to the government and the purchaser. This standard error is usually expressed as a percentage of the total volume in the sale.

*The Economics of Sample Log Scaling.* An administrator attempting to use these tables must first decide what standard error should be selected as desirable. If he demands a very precise estimate the size of the sample becomes so large as to yield no appreciable savings. If a broad limit of error is used there is a risk of considerable underpayment to the government for the timber or, what is more important in actual practice, risk of considerable overcharge to the purchaser. This raises a problem which is familiar in acceptance sampling work. The producer is satisfied to maintain his average quality over a long run, whereas the purchaser wants assurance that the particular lot he is buying is an acceptable one and not one of the unusual defective lots. In the case of sample log scaling a solution should be found for this problem which satisfies both the government and the purchaser; but before attempting a solution certain conditions and assumptions can be set up which will help to point the way to the answer.

First, consideration must be given to the question of who pays for the scaling, the government or the purchaser of the timber? In reality the government bears the cost of scaling and this cost is not added to the purchase price of the logs. Therefore, since the cost to him is nothing, the purchaser might well demand complete or 100 percent scaling. On the other hand, the government should be unwilling to spend more on scaling than necessary since the scaling serves no purpose other than to permit a financial settlement. Thus a reasonable attitude on the part of the government would be to agree to spend as much on scaling as the purchaser would be justified in spending if he were bearing the cost. In seeking a solution to the problem of how much should be spent on scaling, it is logical therefore to use the

assumption that the purchaser is paying for the scaling.

Second, the frequency with which purchases are made will also affect the solution. From the government's point of view a large number of sales is involved and so long as the government receives its proper total payment over all sales the major purpose of scaling is achieved. However, this would certainly not satisfy the purchasers, many of whom may buy timber from the government very infrequently. Furthermore, since certain contributions to state and county governments are made from the timber sale receipts and the size of these contributions is based on the volume of timber cut in local areas, reasonable accuracy for individual sales is necessary. Therefore, in order to protect infrequent purchasers and local governmental units it appears logical to assume that the problem must be solved for a single sale.

If these fundamental assumptions are made the problem can now be stated in the form of a direct question—How much money should a purchaser of timber spend for scaling, assuming that he is making a single purchase? His purpose in scaling the timber is primarily to avoid the risk of over-payment and if we assume that he lacks any urge to gamble on the probability that a sample scale will be less than the true scale then obviously he will want to reduce the risk of over-payment as much as is possible economically. The more he spends for scaling the smaller will be the probable size of over-payment which he is risking. But each dollar's worth of decrease in the range of possible over-payment costs him more than the preceding dollar's worth of decrease. Therefore, a point will be reached at which he is spending more for additional scaling than he can possibly save by further reduction in the risk of over-payment.

The first step in solving the problem is to express the limit of error in terms of dollars rather than percentage or board feet; and the next step is to convert the number of sample logs which must be scaled to give this limit of error, into the total cost of scaling that number of logs. The following notation is used.

$V$ = the coefficient of variation.

$C$ = the cost of scaling per log in dollars.

$M$ = approximate total value of logs in sale in dollars.

$T$ = total spent for scaling in dollars.

$n$ = number of logs in sample.

$X$ = limit of error of total sale in dollars.

For practical purposes we will use a limit of error that is twice the standard error.

The total cost T, of scaling, will be:

$$T = Cn, \tag{1}$$

whence $n = T/C,$ \hfill (1a)

and $X = 2VM/\sqrt{n}.$ \hfill (2)

Substituting T/C for n in (2) we get

$$X = 2VM\sqrt{C/T} \tag{3}$$

Formula (3) shows the largest error (X) likely to be made when a given sum (T) is spent for scaling. As the total amount spent for scaling increases, the limit of error decreases. Any increase in T will produce a decrease in X. Obviously there is no point in adding a dollar to T, the cost of scaling, unless the reduction in the limit of error (X) is at least a dollar. The point where scaling should be stopped, therefore, is the point at which for a small increase in total spent for scaling (T) there is an equal decrease in the limit of error (X). This point can be found by setting the derivative of X with respect to T equal to —1 and solving for T as in formulae (4) and (5).

$$T = C\sqrt[3]{(VM/C)^2}. \tag{4}$$

Substituting Cn for T and dividing both sides by C we get

$$n = \sqrt[3]{(VM/C)^2}. \tag{5}$$

If more logs are scaled, then money is being wasted because the additional cost of the scaling is not equaled by the reduction in the possible error in payment. If fewer logs are scaled then advantage is not being taken of the possibility of reducing the error at a cost less than the reduction achieved.

For practical application, tables can be drawn up based on the solution of equation (5) which show the number of logs it is reasonable to scale for different combinations

of coefficient of variation, total value of sale, and unit cost of scaling.

*Correction for Finite Population.* In deriving formula (5) for the number of logs to be scaled no account was taken of the fact that the logs are drawn from a finite population without replacement. The correction for sampling without replacement can easily be applied to the expression for X in (2), giving

$$X = 2VM\sqrt{(N-n)/Nn}, \qquad (6)$$

where $N$ = approximate total number of logs in the sale.

As before, substitution of T/C for n gives

$$X = 2VM\sqrt{(NC-T)/NT}. \qquad (7)$$

The derivative of this expression, however, is less simple, and when it is set equal to —1, the following quartic equation results:

$$CNT^3 - T^4 = N(MVC)^2. \qquad (8)$$

For practical purposes this equation may be solved approximately for T.

The first approximation gives

$$T_1 = C\sqrt[3]{(VM/C)^2}. \qquad (9)$$

It will be seen that this is exactly equal to the formula for T when the finite correction factor is not taken into account. A second approximation to the solution for T gives

$$T_2 = T_1\sqrt[3]{1+n_1/N}. \qquad (10)$$

This equation can be applied still more directly by dividing both sides by C, and since $T_1/C = n_1$, we get

$$n_2 = n_1\sqrt[3]{1+n_1/N}. \qquad (11)$$

The radical is just the correction introduced by the finite correction factor. This factor indicates that it is reasonable to scale a larger sample of logs than is indicated when the correction factor is not used. However, the effect of the correction factor is very small until the size of the sample approaches 50 percent or larger. From the standpoint of theoretical interest, however, the following table is given:

| Percentage of Logs to be Sampled Assuming Unlimited Population | Percentage of Logs to be Sampled Assuming Limited Population |
|---|---|
| 10 | 10.3 |
| 20 | 21.2 |
| 30 | 32.7 |
| 40 | 44.7 |
| 50 | 57.2 |
| 75 | 90.4 |

The first column shows the percentage of the total number of logs which formula (5) indicates should be sampled on the assumption of an unlimited population. The second column shows the corresponding percentage if formula (11), assuming a limited population, were used. From a practical standpoint it is clear that the assumption of an unlimited population is entirely satisfactory because the differences are quite small until the samples become such a large proportion of the total that it is likely that sampling would not represent a saving over 100 percent scaling.

*An alternate method.* A second closely related solution can be found by assuming that it is desired to control the *average* amount of over-payment rather than the *maximum* over-payment which might occur. This solution would appear to be the better one in a case where the purchaser expects to make a number of purchases and is interested in keeping the sum of the possible over-payments on all purchases within a reasonable limit.

For the solution of this case it is necessary to know the number of purchases to be made by each individual. The more purchases he makes the less scaling will be necessary for each one to insure any desired limit to the sum of the over-payments. Since the number of sales to be made usually will not be known some sort of an approximate solution will be necessary. We might assume that the purchaser is concerned only with the gross sum of his over-payments, not the net balance between under- and over-payments. A possible solution in this case would be to minimize the sum of the over-payments plus the cost of scaling. Since the average over-payment is equal to .80V all that is necessary is to substitute .80V for 2V in formula (1) and work through the same type of solution. This gives the number of logs to be scaled equal to

about half the number that would be scaled if the first method were used.

*Conclusion.* The problem of determining the number of logs to be scaled on national forest timber scales is a particular case of the more general problem of determining economic sampling intensities in acceptance inspection sampling. In this case certain postulates have been set up which help in reaching a solution. The first of these is that the purchaser is to pay for the sampling; the second that only a single transaction is involved. Keeping these conditions in mind a reasonable solution would seem to be to sample to the point where the marginal reduction in the limit of error (taken as two standard errors) expressed in monetary units equals the marginal increase in the cost of sampling.

The same procedure can be applied to other problems where a value per unit can be ascribed to the units in which the standard error is expressed.

# QUERIES

**QUERY** The usual rule in analysis of variance is to divide the treatment mean square by the appropriate error, but to use the F table the larger mean square must be divided by the smaller. What do you do if the treatment mean square is the smaller?

**ANSWER** In the majority of experiments to whose results analysis of variance is applied, the investigator is trying to learn if the treatments are effective in differentiating the means of the sampled populations. If the treatments are effective, then the corresponding mean square is greater than that for error, and the rule for entering the table is suitable for testing significance. If F is larger than the tabled 5% (or 1%) one concludes that the treatments have real effects. If, on the other hand, F is smaller, the experiment is usually thought to lack adequate evidence of population differences.

Now, if the treatment mean square is no greater than error, it is clear that F cannot reach the 5% level no matter how large $n_1$ or $n_2$ may be. It is a tenable conclusion, then, that the treatment means exhibit only sampling fluctuation about some common population mean. There is no occasion for going through the motions of making a test because non-significance is obvious.

This method of testing is adapted to the needs of the majority—those who wish to know if some treatments or varieties are more desirable than others: the null hypothesis is rejected only if F falls in the upper 5% of the distribution. If querist has other alternatives to be considered, he may be interested in the two-tailed test discussed in following answers.

George W. Snedecor.

**QUERY** In the April number of *Industrial and Engineering Chemistry*, Analytic Edition, Mandel was testing the significance of the difference between two variances. He said that in this case a 5% level actually changes into 10% since one half of all possible values of F have been eliminated. What is meant by that statement?

**ANSWER** Mandel was testing the hypothesis that the variances of two independent samples are estimates of a common $\sigma^2$. Since in this case one must wait until the experimental results are determined before learning which variance is the larger, he must allow for rejection of the hypothesis if F falls in either tail of the distribution. This is done by doubling the tabular probabilities corresponding to F calculated in the usual way—divide the larger variance by the smaller. If Mandel had wished to amplify his statement, he might have said something like this: There is no prior information about the variances, hence the region of rejection of the null hypothesis must include both large and small values of F; only half of

70

these values lie beyond the 5% point found in the table, the eliminated half being in the lower end of the distribution, not tabulated; hence, to arrive at the correct probability, read 10% instead of the tabular 5%.

<div align="right">George W. Snedecor.</div>

**QUERY** I understand that in an F test the two mean squares must be independent. Can I test the ratio of two mean squares which are correlated?

**ANSWER** In certain cases, yes, though the test used is not the customary F test. A test was developed by Pitman and Morgan (*Biometrika* XXXI, 1939, p. 9). The following example illustrates its application.

The data are the average wheat yields, in bushels per acre, for two Nebraska counties during the 12 years, 1926-1937.

<div align="center">WHEAT YIELDS (BU.)</div>

| Year | Sarpy | Lancaster |
|------|-------|-----------|
| 26 | 18.3 | 16.0 |
| 27 | 19.9 | 21.9 |
| 28 | 20.5 | 17.3 |
| 29 | 18.6 | 20.8 |
| 30 | 20.3 | 19.0 |
| 31 | 20.6 | 22.2 |
| 32 | 21.2 | 13.7 |
| 33 | 19.3 | 18.7 |
| 34 | 12.0 | 9.0 |
| 35 | 16.7 | 16.6 |
| 36 | 20.8 | 20.4 |
| 37 | 20.5 | 18.8 |
| Mean | 19.1 | 17.9 |

It is desired to test whether the variability in yield from year to year is the same in the two counties. Since the counties are near one another and thus subject to similar weather conditions, a correlation between their annual yields is to be expected. This will produce a correlation between the two mean squares and vitiate the conditions for the ordinary F test.

In Pitman and Morgan's test, we start with the assumption that the two sets of yields follow a bivariate normal distribution, the size of the true correlation coefficient being unknown. The sums of squares and products of deviations are

$$SS \ (Sar.) = 72.63,$$
$$SS \ (Lan.) = 154.31,$$
$$SP = 73.34,$$

Thus the sample correlation r between the the yields is found to be .6925, which with 12 pairs is significant at the 2 percent level.

The value of F is $154.31/72.63$, or 2.125, but instead of referring this to the usual table, we calculate the quantity

$$(F-1)/\sqrt{(F+1)^2 - 4r^2F}$$

$$1.125/\sqrt{(3.125^2 - 4(.6925)^2 \cdot 2.125} = .472.$$

Rather curiously, this quantity is distributed as a sample correlation coefficient from 12 pairs of observations, and the variance-ratio test is made by seeing whether $r = .472$ is significant from 12 pairs.

One further point must be considered. If, so far as we know beforehand, *either* county may be the more variable, then we consult the 5 percent level of r in the table in order to make a 5 percent significance test. If, on the other hand, we are testing a specified county, chosen before seeing the data, is the more variable, r is referred to the 10 percent level in the correlation coefficient table. In my experience the former situation is the more common in this type of test. In this example the observed r, .472, falls considerably short of the 5 percent level of r in the table. With 12 pairs, the 10 percent level of r is .497. Consequently, if we are making the first type of test, we conclude that the variance ratio just fails to attain significance at the 10 percent level. If we are making a one-sided test, on the assumption that Lancaster County is certain to be at least as variable as Sarpy County, the conclusion is that the variance ratio just fails to reach 5 percent significance.

If r is small, this test is less sensitive than the erroneous F test. This is to be expected, because the F test requires the additional assumption that the true correlation is zero. If, however, r is high, either positively or negatively, the present test is actually more powerful than the erroneous F test.

<div align="right">W. G. Cochran.</div>

# NEWS AND NOTES

EDWIN J. DE BEER. The Wellcome Research Laboratories, says that L. O. RANDALL of their research staff gave him a bit of information which will probably interest our readers. In 1936 a "Biometric Bulletin" was published at the Worcester State Hospital, Worcester, Mass., and dealt with the activities of the Memorial Foundation for Neuro-Endocrine Research and Research Service of the Worcester State Hospital. E. MORTON JELLINEK was the editor and this Bulletin was published quarterly for one year only . . . HOWARD V. JORDAN, Associate Soil Technologist, Bureau of Plant Industry, has recently been stationed at Miss. State College. Since 1942, when he left Texas, he has been at Madison, Wisconsin on a war-crop assignment (hemp) . . . CLARENCE DORMAN, Director of the Agricultural Experiment Station of Miss. State College, has been named acting president . . . S. R. MILES, Agricultural Experiment Station, Purdue Univ., has been appointed as an adviser on design of experiments for the station staff members. That 1941 statistical summer session group might be interested to hear that plans are being formulated for another session during the summer of 1946 . . . On June 18, BESSE DAY spoke at the Buffalo meeting of the Society of Quality Control Engineers . . . PAUL PEACH, Industrial Statistician, Institute of Statistics, Raleigh, North Carolina, is to conduct a course in "Industrial Statistics and Quality Control" October 16-26. His instructional staff includes W. EDWARDS DEMING, Adviser in Sampling, Bureau of the Budget, RALPH HEFNER, Professor of Mathematics, Georgia School of Technology and FREDERICK MOSTELLER, Research Statistician, Statistical Research Group, Columbia University. The students are in for a siege with two master magicians as teachers . . . R. A. FISHER, Department of Genetics, Cambridge University did visit India and Egypt last winter . . . ARTHUR B. CHAPMAN, Assistant Professor of Genetics, University of Wisconsin, says his copies of the Bulletin are already shopworn due to the use which has been made of them by his students. We like that! . . . J. A. HALL, who has served as principal biochemist for the past three years for the U. S. Forest Service, Washington, D. C., has become director of the Pacific Northwest Forest Experiment Station at Portland, Oregon. Those men in the field of forestry have helped pull our Bulletin through the first five issues, by contributing three articles. Some of the rest of you readers take note. We welcome good articles for consideration . . . That traveler, C. F. SARLE has just returned from a six weeks trip to Alaska and the Pacific Coast. He visited airway and district forecasting centers trying to arouse interest in objective methods of forecasting . . . Birthday greetings to G. W. SNEDECOR, the date is October 20. Greet him with some queries!

Officers of the American Statistical Association, President: Walter A. Shewhart, Directors: Henry B. Arthur, C. I. Bliss, Simon Kuznets, E. Grosvenor Plowman, Willard L. Thorp, and Helen M. Walker; Vice-Presidents, William G. Cochran, A. D. H. Kaplan, Lowell J. Reed; Secretary-Treasurer, Lester S. Kellogg.

Officers of the Biometrics Section: C. I. Bliss, Chairman; H. W. Norton, Secretary.

Editorial Committee for the BIOMETRICS BULLETIN: Gertrude M. Cox, Chairman; C. I. Bliss, W. G. Cochran, F. R. Immer, J. Neyman, H. W. Norton, L. J. Reed, G. W. Snedecor, Sewall Wright.

Material for the BULLETIN should be addressed to the Chairman of the Editorial Committee, Institute of Statistics, North Carolina State College, Raleigh, N. C., material for Queries should go to "Queries", Statistical Laboratory, Iowa State College, Ames, Iowa, or to any member of the committee.

B/B